

# NINTH LECTURE SUMMARY

## REGRESSION AND CORRELATION

HERE, WE STUDY RELATIONSHIP BETWEEN TWO VARIABLES (EACH ON AN INTERVAL SCALE) USUALLY CALLED  $X$  (EXPLANATORY) AND  $Y$  (RESPONSE).

### ASSUMPTIONS:

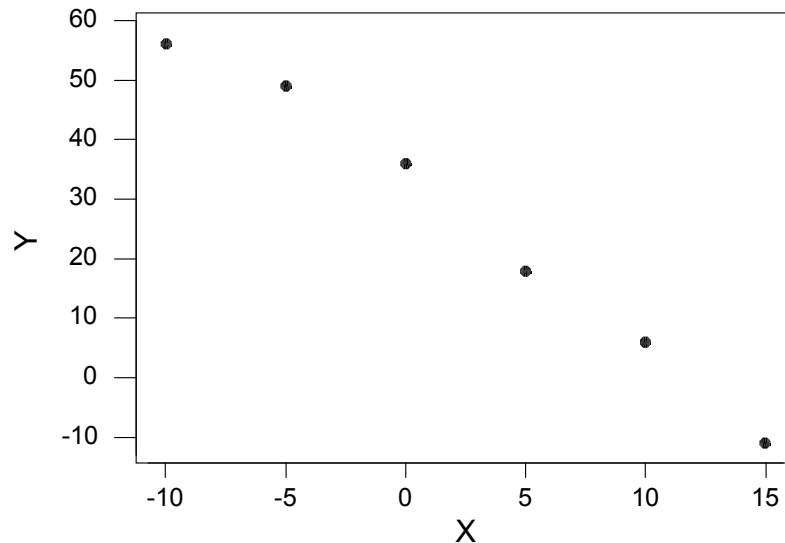
- RELATIONSHIP IS LINEAR (FOLLOWING A STRAIGHT LINE)
- RESIDUALS (RANDOM DEVIATIONS OF INDIVIDUAL  $Y$  VALUES FROM THIS STRAIGHT LINE) ARE NORMALLY DISTRIBUTED

SAMPLE HAS A FORM OF A TABLE, E.G.

$X$ ( $^{\circ}\text{C}$ )	-10	-5	0	5	10	15
$Y$ (\$)	56	49	36	18	6	-11

( $Y$  IS DAILY NET PROFIT OF A TEENAGE ENTREPRENEUR SELLING HOT CHOCOLATE,  $X$  IS THE CORRESPONDING DAY'S AVERAGE TEMPERATURE).

THE DATA CAN BE DISPLAYED GRAPHICALLY IN SO CALLED **SCATTER DIAGRAM**:



TO ANSWER ALL POTENTIAL QUESTIONS, WE FIRST COMPUTE

$E_x$ ,  $E_y$ ,  $E_{x^2}$ ,  $E_{y^2}$ ,  $E_{xy}$

THEN CONVERT THEM TO

$$\bar{x} = \frac{\sum x}{n} \quad \bar{y} = \frac{\sum y}{n} \quad SS_x = \sum x^2 - \frac{(\sum x)^2}{n}$$

$$SS_y = \sum y^2 - \frac{(\sum y)^2}{n} \qquad SP_{xy} = \sum x \cdot y - \frac{(\sum x) \cdot (\sum y)}{n}$$

HAVING COMPUTED THESE FIVE BASIC QUANTITIES, WE CAN NOW FIND THE **REGRESSION COEFFICIENTS** (INTERCEPT AND SLOPE) OF THE BEST (**LEAST SQUARES**, REGRESSION) STRAIGHT LINE:

$$b = \frac{SP_{xy}}{SS_x} \quad (\text{SLOPE FIRST}) \qquad a = \bar{y} - b \cdot \bar{x}$$

NOTE THAT THIS LINE MUST PASS THROUGH THE  $(\bar{x}, \bar{y})$  POINT.

NEXT, WE NEED TO COMPUTE **RESIDUAL STANDARD DEVIATION** (TYPICAL MAGNITUDE OF THE RESIDUALS) BY:

$$s_r = \sqrt{\frac{SS_y - b \cdot SP_{xy}}{n - 2}}$$

WHERE **RESIDUALS** ARE DEFINED AS

$$y_i - (a + b \cdot x_i)$$

THE **PREDICTION INTERVAL** (ALSO CALLED **CONFIDENCE INTERVAL FOR PREDICTION**) FOR A NEW **Y** OBSERVATION, TAKEN AT **X** =  $x_0$  (A SPECIFIC NUMBER) IS:

$$a + b \cdot x_0 \pm t_c \cdot s_r \cdot \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_x}}$$

WHERE  $t_c$  IS THE CORRESPONDING CRITICAL VALUE OF THE STUDENT DISTRIBUTION, USING  $n - 2$  DEGREES OF FREEDOM.

SAMPLE **CORRELATION COEFFICIENT** IS COMPUTED BY

$$r \equiv \frac{SP_{xy}}{\sqrt{SS_x \cdot SS_y}}$$

(ALWAYS BETWEEN -1 AND 1).

COEFFICIENT OF DETERMINATION:  $r^2$   
(PERCENTAGE REDUCTION OF THE ORIGINAL  $y$  VARIANCE).

## TWO TESTS

$H_0: \beta = 0$  ..... WHERE  $\beta$  IS THE SLOPE OF THE POPULATION STRAIGHT LINE

TEST STATISTIC:  $\frac{b}{s_r} \cdot \sqrt{SS_x}$  (USE  $t_{n-2}$ )

$H_0: D = 0$  ..... WHERE  $D$  IS THE POPULATION CORRELATION COEFFICIENT

TEST STATISTIC:  $\frac{r \cdot \sqrt{n-2}}{\sqrt{1-r^2}}$  (USE  $t_{n-2}$ )

ONE CAN SHOW THAT THE TWO TEST STATISTICS (THUS THE TWO TESTS) ARE IDENTICAL.